

# Pokročilé statistické metody pro biology

Alena Černíková

alena.cernikova@ujep.cz

4. ledna 2024

- **tři domácí úkoly**

jednoduché opakování příkladů ze cvičení  
odevzdávat na univerzitní OneDrive – bude upřesněno  
později

důraz je kladen na interpretaci výsledků

- **seminární práce**

zpracování několika proměnných  
od zadání *výzkumu* až po interpretaci

- Opakování – popisné statistiky, bodové a intervalové odhady, jednovýběrové, dvouvýběrové a párové t-testy
- Neparametrické testy
- Analýza rozptylu
- Korelační koeficienty
- Regresní modely – mnohonásobná regrese, zobecněná lineární regrese
- Kontingenční tabulky – poměr šancí, testy pro ordinální data
- Mnohorozměrné statistické metody

- **Nahodná veličina** – jakákoliv veličina, kterou měříme, např. výška
- **Populace** – soubor, pro nějž chceme udělat nějaký závěr, např. všichni dospělí obyvatelé České republiky
- **Náhodný výběr** – v porovnání s populací malý soubor pozorování, jde o nezávislé, stejně rozdělené náhodné veličiny, např. výběr 200 lidí
- **Statistická jednotka** – objekt, na kterém měříme, např. člověk
- **Populační charakteristika** – charakteristika popisující populaci, např. populační průměr
- **Výberová charakteristika** – charakteristika spočítaná na výběru pomocí níž odhadujeme populační ekvivalent, např. výběrový průměr.

## ● Číselné proměnné

- popisné statistiky polohy – průměr, medián, vybrané percentily (kvartily, extrémy)
- popisné statistiky variability – rozptyl, směrodatná odchylka, mezikvartilové rozpětí, koeficient variace
- popisné statistiky tvaru rozdělení – šikmost, špičatost
- grafické charakteristiky – krabicový graf, histogram

## ● Nominální proměnné

- číselné charakteristiky – absolutní a relativní četnosti
- grafické charakteristiky – sloupcový a koláčový graf

## ● Ordinální proměnné

- lze použít jak průměr, medián atd.
- a pro malé počty kategorií i absolutní a relativní četnosti

Popisné statistiky polohy pro číselnou proměnnou

- **průměr** –  $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ , kde  $n$  je počet pozorování a  $X_1, X_2, X_3, \dots, X_n$  jsou jednotlivá měření
- **medián** – hodnota prostřední podle velikosti, nebo průměr prostředních dvou
- vybrané percentily, především **kvartily** – hodnoty v jedné a ve třech čtvrtinách podle velikosti

$$p\text{-ty percentil} = (1 - q)X_{(k)} + qX_{(k+1)}$$

$$k = \lfloor 1 + (n - 1)p \rfloor, q = 1 + (n - 1)p - k$$

Popisné statistiky variability pro číselnou proměnnou

- **Rozptyl** –  $\text{Var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$
- **Směrodatná odchylka** –  $\text{sd}(X) = \sqrt{\text{Var}X}$
- **Mezikvartilové rozpětí** –  $\text{IQR}(X) = Q_3 - Q_1$ , kde  $Q_3$  je třetí kvartil a  $Q_1$  je první kvartil
- **Variační koeficient** –  $\text{cv}(X) = \frac{\text{sd}(X)}{\bar{X}}$
- **Rozpětí** –  $\max(X) - \min(X)$
- **Střední absolutní odchylka** –  $\text{MAE}(X) = \frac{\sum_{i=1}^n |X_i - \bar{X}|}{n}$

Popisné statistiky tvaru rozdělení pro číselnou proměnnou se počítají ze standardizovaných proměnných, tak zvaných **Z-skórů**

$$Z_i = \frac{X_i - \bar{X}}{\text{sd}(X)}$$

- **Šikmost** – průměr ze třetích mocnin z-skórů

$$\text{Skew}(X) = \frac{1}{n} \frac{\sum_{i=1}^n (X_i - \bar{X})^3}{\text{sd}(X)} = \frac{\sum_{i=1}^n Z_i^3}{n}$$

- **Špičatost** – průměr ze čtvrtých mocnin z-skórů mínus 3

$$\text{Kurt}(X) = \frac{1}{n} \frac{\sum_{i=1}^n (X_i - \bar{X})^4}{\text{sd}(X)} - 3 = \frac{\sum_{i=1}^n Z_i^4}{n} - 3$$



## Popisné statistiky pro nominální proměnnou

- absolutní četnosti – kolik hodnot se naměřilo v dané kategorii,  $n_i$
- relativní četnosti – udávají se buď v desetinných číslech nebo v procentech  $p_i = n_i/n$

## Senzitivita a specificita testu

- **Senzitivita testu** – pravděpodobnost, že test vyjde pozitivně, pokud je osoba nemocná  
 $P(\text{test je pozitivní} | \text{osoba je nemocná})$
- **Specificita testu** – pravděpodobnost, že test vyjde negativně, pokud je osoba zdravá  
 $P(\text{test je negativní} | \text{osoba je zdravá})$

Pravděpodobnostní rozdělení dělíme podle typu proměnné na

- **Spojité** – pro číselné proměnné, které teoreticky mohou nabývat libovolné reálné hodnoty z nějakého intervalu, př. normální, exponenciální, chí-kvadrát, ...
- **Diskrétní** – pro kategorické proměnné s jasně oddělitelnými kategoriemi, může být i nekonečně mnoho hodnot  
př. binomické, poissonovo, alternativní, ...

## Funkce určující rozdělení

- **Distribuční funkce** –  $F(t) = P(X \leq t), t \in \mathbb{R}$ 
  - neklesající, zprava spojitá, obor hodnot je mezi 0 a 1
- **Pravděpodobnostní funkce** –  $p(t) = P(X = t), t \in \mathbb{R}$ 
  - definovaná pouze pro diskrétní rozdělení
  - nespojitá, nenulová jen v hodnotách, kterých může náhodná veličina nabývat
- **Hustota** –  $f(t) = \frac{d}{dt}F(t)$ 
  - definovaná pouze pro spojitá rozdělení – obdoba pravděpodobnostní funkce, ale nedefinuje konkrétní pravděpodobnosti
  - derivace funkce distribuční
  - pravděpodobnost jedné konkrétní hodnoty u spojitého rozdělení je 0

# Pravděpodobnostní rozdělení

## Binomické rozdělení – zástupce diskretních rozdělení

Mějme náhodný pokus, který může skončit jedním ze dvou výsledků: úspěch – neúspěch. Opakujme tento pokus mnohokrát a počítejme počet úspěchů. Počet úspěchů má binomické rozdělení.

Značení  $Bi(n, p)$ , kde

- $n$  – počet pokusů,
- $p$  – pravděpodobnost úspěchu

Hodnoty pravděpodobnostní funkce

$$p(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Střední hodnota a rozptyl

$$E(X) = np,$$

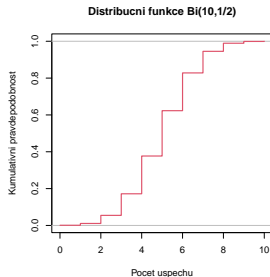
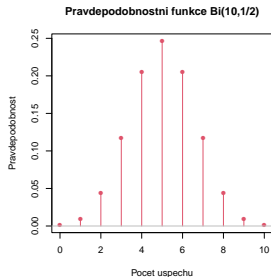
$$\text{Var}(X) = np(p - 1)$$

# Pravděpodobnostní rozdělení

## Binomické rozdělení

**Příklad.** *Házíme 10x mincí a počítáme, kolikrát padla panna. Počet pokusů je  $n = 10$ , pravděpodobnost úspěchu  $p = 1/2$ . Máme tedy rozdělení  $Bi(10, 1/2)$ .*

Pravděpodobnostní a distribuční funkce.



Střední hodnota a rozptyl

$$E(X) = np = 10 \frac{1}{2} = 5, \quad \text{Var}(X) = np(1 - p) = 10 \frac{1}{2} \frac{1}{2} = 2.5$$

# Pravděpodobnostní rozdělení

## Poissonovo rozdělení – zástupce diskretních rozdělení

**Př.** Sledujeme počet nehod na křižovatce v průběhu jednoho dne. Za normálních okolností nenastane ani jedna nehoda, nebo nastane jedna, maximálně 2 nehody. Ale může se stát, že při náledí jich nastane klidně i 10. Tato veličina má Poissonovo rozdělení

Značení  $Po(\lambda)$ , kde

- $\lambda$  – parametr rozdělení, intenzita

Hodnoty pravděpodobnostní funkce pro  $k = 0, 1, 2, \dots$

$$p(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}$$

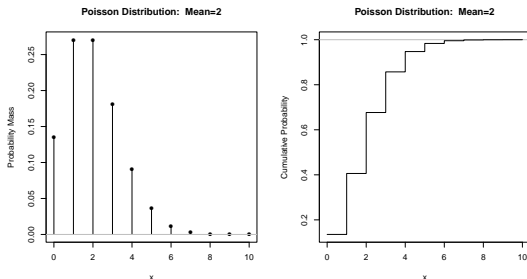
Střední hodnota a rozptyl

$$E(X) = \lambda,$$

$$\text{Var}(X) = \lambda$$

## Poissonovo rozdělení

Pravděpodobnostní a distribuční funkce Poissonova rozdělení s parametrem  $\lambda = 2$ .



Předpokládejme binomická rozdělení  $Bi(n, p_n)$ , kde  $np_n \rightarrow \lambda$ , pak tato binomická rozdělení konvergují k rozdělení Poissonovu s parametrem  $\lambda$

# Pravděpodobnostní rozdělení

## Hypergeometrické rozdělení – zástupce diskrétních rozdělení

**Př.** Uvažujme urnu, ve které máme  $N$  koulí, z toho  $A$  jich je bílých a zbytek černých. Z urny postupně vytáhneme  $n$  koulí bez vracení. Náhodná veličina, která počítá počet bílých koulí mezi vytaženými má hypergeometrické rozdělení.

Značení  $Hy(N, A, n)$ , kde

- $N$  – počet koulí v urně
- $A$  – počet označených koulí v urně
- $n$  – počet tažených koulí

Hodnoty pravděpodobnostní funkce

$$p(X = k) = \frac{\binom{A}{k} \binom{N-A}{n-k}}{\binom{N}{n}}$$

Střední hodnota a rozptyl

$$E(X) = \frac{nA}{N},$$

$$\text{Var}(X) = n \frac{A}{N} \left(1 - \frac{A}{N}\right) \left(\frac{N-n}{N-1}\right)$$



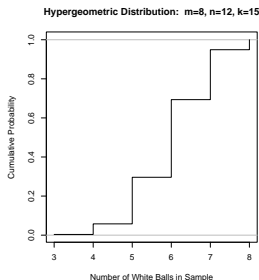
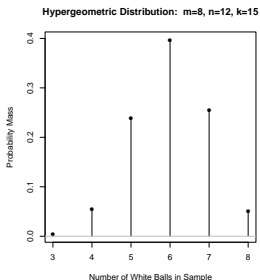


# Pravděpodobnostní rozdělení

## Hypergeometrické rozdělení

**Příklad.** Uvažujme 20 koulí v urně, z toho 8 bílých a z urny vytáhneme 15 koulí .

Pravděpodobnostní a distribuční funkce.



Střední hodnota a rozptyl

$$E(X) = \frac{nA}{N} = 6, \quad \text{Var}(X) = n \frac{A}{N} \left(1 - \frac{A}{N}\right) \left(\frac{N-n}{N-1}\right) = 0.95$$

## Normální rozdělení – zástupce spojitých rozdělení

Jedná se o "hezké" rozdělení, se kterým se dobře pracuje. Toto rozdělení má výška lidí určitého věku, IQ, . . . . Ve statistice se nejčastěji používá standardní normální rozdělení  $N(0, 1)$

Značení  $N(\mu, \sigma^2)$ , kde

- $\mu$  – střední hodnota
- $\sigma^2$  – rozptyl

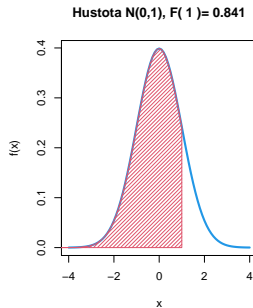
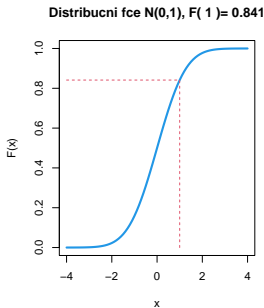
Hustota normálního rozdělení má tvar

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x - \mu)^2}{2\sigma^2}\right\}$$

Je to tak zvaná **Gaussova křivka**.

## Normální rozdělení

Vztah mezi hustotou a distribuční funkcí u standardního normálního rozdělení  $N(0, 1)$ . Červeně je na obou grafech zobrazena stejná hodnota. Hustota a distribuční funkce.



Předpokládejme binomické rozdělení  $Bi(n, p)$ , kde  $0.1 \leq p \leq 0.9$ , pak pro  $n \rightarrow \infty$  toto rozdělení konverguje k normálnímu s parametry  $np, np(1 - p)$ .

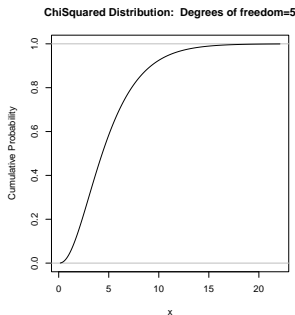
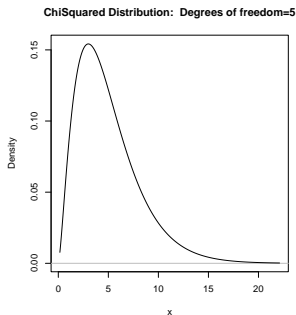
# Pravděpodobnostní rozdělení

$\chi^2$ -rozdělení – zástupce spojitych rozdělení

Rozdělení kvadratických forem. Náhodná veličina

$Y = X_1^2 + X_2^2 + \dots + X_n^2$ , kde  $X_i \sim N(0, 1)$  jsou nezávislé, má

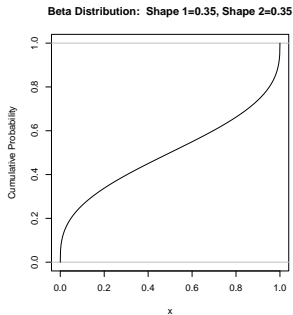
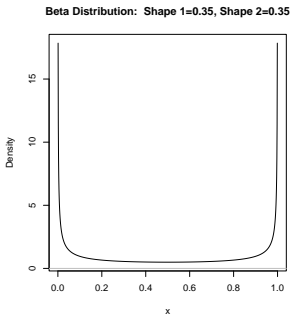
$\chi^2$ -rozdělení o  $n$  stupních volnosti. Dále je to rozdělení některých testových statistik, zejména těch, týkajících se rozptylu. Hustota a distribuční funkce  $\chi^2$ -rozdělení o 5 stupních volnosti



# Pravděpodobnostní rozdělení

## Beta – zástupce spojitých rozdělení

Rozdělení pravděpodobností nějakého jevu. Např. sledujeme pravděpodobnost, že vybrný člověk má nebo nemá nějakou nemoc. Rozdělení má 2 tvarové parametry, které určují, jak vypadají pravděpodobnosti u 0 a 1. Hustota a distribuční funkce Beta rozdělení s parametry 0.35 a 0.35



Statistika se zabývá odhadem/testováním teoretických/populačních charakteristik. Nejčastěji odhadujeme

- **pravděpodobnost** náhodného jevu –  $\pi_i$
- **Střední hodnotu** –  $E(X) = \sum_{i=1}^n X_i p_i = \int_{-\infty}^{\infty} x f(x) dx$   
– vlastnosti:  $E(aX + b) = aE(X) + b$ ,  $E(X + Y) = E(X) + E(Y)$
- **Rozptyl** –  $\text{Var}(X) = \sum_{i=1}^n (X_i - E(X))^2 p_i = \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx$   
– vlastnosti:  
 $\text{Var}(aX + b) = a^2 \text{Var}(X)$ ,  $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y)$
- **Korelace** –  $\text{cor}(X, Y) = \text{cov}(X, Y) / (\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}) =$   
 $\frac{\sum_{i=1}^n (X_i - E(X))(Y_i - E(Y)) p_i q_i}{(\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)})} =$   
 $\frac{\int_{-\infty}^{\infty} (x - E(X))(y - E(Y)) f(x, y) dx dy}{(\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)})}$

kde  $X, Y$  jsou náhodné veličiny,  $X_i, Y_i$  hodnoty, jichž mohou nabývat u diskrétní veličiny,  $p_i, q_i$  pravděpodobnosti hodnot,  $f(\cdot)$  je hustota

## Odhad pravděpodobnosti

- nejlepším bodovým odhadem pravděpodobnosti je relativní četnost  $p_i = n_i/n$
- nestranný odhad
- náhodná veličina  $p = (p_i - \pi_i)/\sqrt{\pi_i(1 - \pi_i)/n}$  konverguje k normálnímu rozdělení  $N(0, 1)$  pro  $n \rightarrow \infty$
- intervalový odhad pro pravděpodobnost je

$$\left( p_i - z(1 - \alpha/2)\sqrt{p_i(1 - p_i)/n}, p_i + z(1 - \alpha/2)\sqrt{p_i(1 - p_i)/n} \right)$$

- pro použití tohoto intervalu musíme mít dostatečně velké  $n$  a  $p_i$ , má platit  $np_i(1 - p_i) > 9$

Odhad **střední hodnoty**/ populačního průměru

- nejlepším bodovým odhadem střední hodnoty je výběrový průměr  $\bar{X} = \sum_{i=1}^n X_i/n$
- nestranný odhad
- platí **Centrální limitní věta** – pro rostoucí počet pozorování konverguje rozdělení výběrového průměru k normálnímu pro  $n \rightarrow \infty$
- střední chyba průměru je  $SEM = sd(X)/\sqrt{n}$
- intervalový odhad pro průměr je

$$(\bar{X} - t_{n-1}(1 - \alpha/2)sd(X)/\sqrt{n}, \bar{X} + t_{n-1}(1 - \alpha/2)sd(X)/\sqrt{n})$$



Odhad **teoretického rozptylu**/ populačního rozptylu

- jako bodový odhad populačního rozptylu používáme výběrový rozptyl  $\text{Var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$
- nestranný odhad
- označme výběrový rozptyl jako  $s^2$  a teoretický rozptyl jako  $\sigma^2$ , pak náhodná veličina  $\chi = (n-1)s^2/\sigma^2$  má  $\chi^2$  rozdělení o  $n$  stupních volnosti
- $\chi^2$  rozdělení není symetrické
- intervalový odhad pro rozptyl je

$$\left( \frac{(n-1)s^2}{\chi_{n-1}^2(1-\alpha/2)}, \frac{(n-1)s^2}{\chi_{n-1}^2(\alpha/2)} \right)$$

## Odhad teoretického korelačního koeficientu

- nejlepším bodovým odhadem Pearsonova korelačního koeficientu je výběrový korelační koeficient

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

- máme-li dvourozměrné normální rozdělení a odhadovaný korelační koeficient  $|\rho| < 0.5$  pak je interval spolehlivosti pro korelační koeficient

$$\left( \text{Cor}(X, Y) - z(1 - \alpha/2) \frac{1 - \text{Cor}(X, Y)^2}{\sqrt{n-3}}, \right. \\ \left. \text{Cor}(X, Y) + z(1 - \alpha/2) \frac{1 - \text{Cor}(X, Y)^2}{\sqrt{n-3}} \right)$$

## Odhad teoretického korelačního koeficientu – pokračování

- nejsou-li splněny podmínky výše, pak je intervalový odhad pro korelační koeficient odvozen z faktu, že náhodná veličina

$$Z = \frac{1}{2} \ln \left\{ \frac{1 + \text{Cor}(X, Y)}{1 - \text{Cor}(X, Y)} \right\} \sim N \left( \frac{1}{2} \ln \left\{ \frac{1 + \rho}{1 - \rho} \right\} + \frac{\rho}{2(n-1)}, \frac{1}{n-3} \right)$$

- interval spolehlivosti tedy je

$$\left( \text{tgh}(Z - z(1 - \alpha/2)/\sqrt{n-3}), \text{tgh}(Z + z(1 - \alpha/2)/\sqrt{n-3}) \right)$$

kde  $\text{tgh}(x) = (e^x - e^{-x})/(e^x + e^{-x})$

Určení rozsahu výběru na základě požadované délky intervalu spolehlivosti

Předpokládejme, že chceme realizovat výzkum, jehož cílem je odhadnout střední hodnotu s požadovanou přesností. Délka intervalu spolehlivosti nesmí přesáhnout hodnotu  $2\Delta$ .

Platí

$$\Delta \geq z(1 - \alpha/2) \frac{\sigma}{\sqrt{n}}$$

Rozsah výběru pak musí splňovat

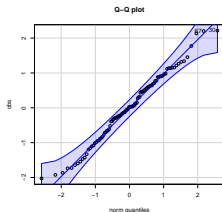
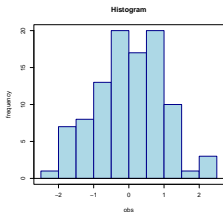
$$n \geq \left( z(1 - \alpha/2) \frac{\sigma}{\Delta} \right)^2$$

# Testování normality

Statistické testy a metody, které znáte, jsou odvozeny pro normální rozdělení. Je proto dobré zjistit, zda náhodná veličina normální rozdělení má či nemá.

K tomuto účelu se využívají

- **Grafické testy** – histogram a pravděpodobnostní graf



- **Číselné testy** – např. Shapiro-Wilkův, Andersonův-Darlingův, Kolmogorovův-Smirnovův, Lillieforsův a další

Nejčastěji používané číselné testy normality

- **Shapiro-Wilkův** – test odpovídající pravděpodobnostnímu grafu  
porovnává, jak si odpovídají teoretické percentily pro normální rozdělení a percentily naměřené pro sledovanou proměnnou
- **Kolmogorovův-Smirnovův** – test je založen na maximálním rozdílu empirické distribuční funkce a distribuční funkce normálního rozdělení
- **Andersonův-Darlingův** – test je založen na váženém průměru druhé mocniny rozdílu empirické distribuční funkce a distribuční funkce normálního rozdělení

# Opakování – testování hypotéz

Při statistickém rozhodování testujeme proti sobě 2 hypotézy

- **Nulovou hypotézu** – značíme  $H_0$   
– je v ní vždy pouze jedna varianta
- **Alternativní hypotézu** – značíme  $H_A$   
– obsahuje více možností (např. interval)

Na základě testu uděláme jedno ze dvou rozhodnutí

- Zamítneme nulovou hypotézu – platí alternativa
- Nezamítneme nulovou hypotézu

Při rozhodování můžeme udělat chybu

- chyba prvního druhu – zamítneme  $H_0$ , přestože platí  
– značí se  $\alpha$ , a jmenuje se **hladina významnosti**  
– závažnější z obou chyb
- chyba druhého druhu – nezamítneme  $H_0$ , přestože neplatí  
– značí se  $\beta$  a hodnota  $1 - \beta$  se nazývá **síla testu**  
– za dané hladiny významnosti chceme test co nejsilnější

Testovat mohou buď přes porovnání **testové statistiky** a **kritické hodnoty** (kvantil vybraného teoretického rozdělení), nebo přes porovnání  **$p$ -hodnoty** a **hladiny významnosti**.

Výsledkem testu v počítači je  **$p$ -hodnota**

- aktuální dosažená hladina testu
- pravděpodobnost, že za platnosti  $H_0$  nastal výsledek, jaký nastal, nebo jakýkoliv jiný, který ještě více odpovídá alternativě
- $p$ -hodnota  $\leq \alpha$  potom **ZAMÍTÁME  $H_0$**
- $p$ -hodnota  $> \alpha$  potom **NEZAMÍTÁME  $H_0$**



# Jednovýběrový t-test

Nejjednodušším testem je **jednovýběrový test o střední hodnotě**. Testujeme

- $H_0$  : střední hodnota =  $\mu_0$
- $H_1$  : střední hodnota  $\neq \mu_0$ , nebo  $< \mu_0$ , nebo  $> \mu_0$

Není-li řečeno jinak, testujeme na hladině významnosti  $\alpha = 0.05$  **Testová statistika** jednovýběrového t-testu je

$$T = \frac{\bar{X} - \mu_0}{\text{sd}(X)} \sqrt{n}$$

a za platnosti nulové hypotézy má tato statistika  $t$ -rozdělení o  $n - 1$  stupních volnosti.

Předpokladem jednovýběrového t-testu je, že průměr testované veličiny má normální rozdělení (díky CLV většinou splněno).

# Párový t-test

V případě, že porovnáváme dva závislé výěry, tedy taková data, která tvoří přirozené páry, používá se **párový test**.

Testované hypotézy v něm jsou

- $H_0$  : střední hodnota rozdílu párů  $= \mu_0$
- $H_1$  : střední hodnota rozdílu  $\neq \mu_0$ , nebo  $< \mu_0$ , nebo  $> \mu_0$

Postup testu je takový, že v prvním kroku spočítám rozdíly mezi všemi páry

$$R_i = X_i - Y_i$$

kde  $X_i$  a  $Y_i$  jsou párová měření, a ve druhém kroku se testuje střední hodnota/ průměr tohoto rozdílu běžným **jednovýběrovým testem**.

**Příklad.** *Porovnávám věk otce a matky, srovnávám sílu pravé a levé ruky, srovnávám měření před a po podání nějakého léku, atd.*

# Dvouvýběrový t-test

Pokud porovnávám dva nezávislé výběry (pozorování nemohu napárovat), pak je potřeba použít **dvouvýběrový test**.

Testujeme

- $H_0$  : rozdíl středních hodnot =  $\mu_0$
- $H_1$  : rozdíl středních hodnot  $\neq \mu_0$ , nebo  $< \mu_0$ , nebo  $> \mu_0$

**Testová statistika** dvouvýběrového t-testu pro shodné rozptyly je

$$T = \frac{\bar{X} - \bar{Y} - \mu_0}{S} \sqrt{\frac{n_1 n_2}{n_1 + n_2}}$$

kde

$$S = \frac{1}{n_1 + n_2 - 2} \left( \sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2 \right)$$

a  $n_1, n_2$  je rozsah výběru  $X$ , respektive  $Y$ . Za platnosti nulové hypotézy má tato statistika  $t$ -rozdělení o  $n_1 + n_2 - 2$  stupních volnosti.

# Dvouvýběrový t-test

V případě, že výběry shodné rozdělení nemají, používá se **Welchova varianta dvouvýběrového t-testu**. Její **testová statistika** má tvar

$$T = \frac{\bar{X} - \bar{Y} - \mu_0}{\sqrt{\frac{\text{Var}(X)}{n_1} + \frac{\text{Var}(Y)}{n_2}}}$$

Tato statistika má za platnosti nulové hypotézy  $t$ -rozdělení o  $\nu$  stupních volnosti, kde

$$\nu = \frac{(\text{Var}(X)/n_1 + \text{Var}(Y)/n_2)^2}{\frac{(\text{Var}(X)/n_1)^2}{n_1-1} + \frac{(\text{Var}(Y)/n_2)^2}{n_2-1}}.$$

kritické hodnoty je možno odvodit, přestože  $\nu$  není celé číslo.

# Test shody rozptylů ve dvou výběrech

Chceme-li rozhodnout, kterou variantu dvouvýběrového t-testu máme použít, je nutné zjistit, zda jsou v obou výběrech stejné rozptyly.

Testujeme

- $H_0$  : rozptyly jsou shodné
- $H_1$  : rozptyly se liší

Testová statistika **F-testu pro dva rozptyly** má tvar

$$F = \frac{\text{Var}(X)}{\text{Var}(Y)}$$

a za platnosti nulové hypotézy má  $F$ -rozdělení o  $n_1 - 1$  a  $n_2 - 1$  stupních volnosti.

Porovnáváme-li střední hodnotu ve více než dvou nezávislých výběrech, používá se **analýza rozptylu**. Testujeme

- $H_0$  : všechny střední hodnoty jsou stejné
- $H_1$  : alespoň jedna střední hodnota se liší

Myšlenka spočívá v porovnání variability **mezi výběry** s variabilitou **v rámci výběrů**.

**Příklad.** *Byla měřena koncentrace mědi v těle ryb.*

*Porovnáváno bylo 5 rybníků, kde z každého byl vyloven vzorek alespoň 10-ti ryb. Liší se od sebe tyto rybníky?*

# Analýza rozptylu – ANOVA

Označme  $X_{ij}$   $i$ -té pozorování z  $j$ -tého výběru,  $\bar{X}_i$  průměr  $i$ -tého výběru,  $\bar{X}_{..}$  celkový průměr všech pozorování,  $n_i$  rozsah  $i$ -tého výběru a  $k$  počet výběrů.

Analýza rozptylu rozkládá celkovou variabilitu

$$SST = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{..})^2$$

na variabilitu vysvětlenou výběry (mezi výběry)  $SS_A$  a variabilitu nevysvětlenou (zbytkovou, v rámci výběrů)  $SS_e$ . Platí

$$\begin{aligned} SST &= \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{..})^2 = \\ &= \sum_{i=1}^k n_i (\bar{X}_i - \bar{X}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 = \\ &= SSA + SSe \end{aligned}$$

Výstupem z analýzy rozptylu je tzv. **tabulka analýzy rozptylu**

	Součty čtverců	Stupně volnosti	Průměrné čtverce	Testová statistika	$p$ -hodnota
Faktor $A$	$SSA$	$df_A = k - 1$	$MSA$	$F = MSA/MSe$	$p$
Chyba $e$	$SSe$	$dfe = n - k$	$MSe$		
Celkem	$SST$	$dft = n - 1$			

Za platnosti nulové hypotézy má testová statistika  $F$ -rozdělení o  $k - 1$  a  $n - k$  stupních volnosti.



# Bartlettův test

Předpokladem analýzy rozptylu je shoda rozptylů ve všech výběrech. Tento předpoklad můžeme zkontrolovat např. prostřednictvím **Bartlettova testu**.

Testujeme

- $H_0$  : rozptyly jsou shodné
- $H_1$  : rozptyly se liší

Testová statistika je založena na výběrových rozptylech v každém výběru zvlášť. Označme  $\text{Var}(X)_i$  výběrový rozptyl v  $i$ -tém výběru a

$$S^2 = \frac{\sum_{i=1}^k (n_i - 1) \text{Var}(X)_i}{n - k},$$
$$C = 1 + \frac{1}{3(k-1)} \left( \sum_{i=1}^k \frac{1}{n_i - 1} - \frac{1}{n - k} \right)$$

Testová statistika

$$B = \frac{1}{C} \left( (n - k) \ln S^2 - \sum_{i=1}^k (n_i - 1) \ln \text{Var}(X)_i \right)$$

ta má za platnosti nulové hypotézy  $\chi^2$ -rozdělení o  $k - 1$  stupních volnosti.

# Párové srovnání

Zajímá-li nás, které konkrétní dvojice výběrů se od sebe významně liší, nelze toto zjistit větším počtem běžných dvouvýběrových testů, neboť by tím příliš vzrostla chyba prvního druhu (tj. neudržela by se celková hladina významnosti). Je nutné použít párové srovnání, např. **Tukeyův test**, případně **Tukey HSD test** pro různě velké výběry.

Testuje se

- $H_0$  : střední hodnoty  $\mu_i$  a  $\mu_j$  jsou stejné
- $H_1$  : střední hodnoty  $\mu_i$  a  $\mu_j$  se liší

pro všechny dvojice  $i$  a  $j$ .

Testová statistika má tvar

$$Q = \frac{|\bar{X}_{i.} - \bar{X}_{j.}|}{s^*}, \text{ kde } s^* = \sqrt{\frac{SSe}{n(n-k)}}$$

Rozdělení těchto statistik se jmenuje studentizované rozpětí a má své vlastní tabelované kritické hodnoty.

Nemají-li všechny porovnávané skupiny stejné rozptyly, používá se tzv. **Welchova ANOVA**. Ta je založena na myšlence vážení skupinových průměrů vahou odpovídající jejich variabilitě. Namísto celkového průměru se pracuje s váženým průměrem

$$\bar{X}_w = \frac{\sum_{i=1}^k w_i \bar{X}_i}{\sum_{i=1}^k w_i}, \text{ kde } w_i = \frac{n_i}{\text{Var}(X)_i}$$

Variabilita mezi výběry se pak počítá jako

$$SSA_w = \sum_{i=1}^k w_i (\bar{X}_i - \bar{Y}_w)^2, \quad MSA_w = \frac{SSA_w}{k-1}$$

Dále se zavádí parametr

$$\Lambda = \frac{3 \sum_{i=1}^k \frac{\left(1 - \frac{w_i}{\sum_{i=1}^k w_i}\right)^2}{n_i - 1}}{k^2 - 1}$$

testová statistika pak má tvar

$$F_w = \frac{SSA_w / (k - 1)}{1 + \frac{2\Lambda(r-2)}{3}}$$

kteřá má za platnosti  $H_0$   $F$ -rozdělení o  $r - 1$  a  $1/\Lambda$  stupních volnosti.

# Pearsonův korelační koeficient

Je-li cílem výzkumu zjistit, zda spolu lineárně souvisí dvě číselné proměnné, používá se **korelační koeficient**.

**Pearsonův korelační koeficient** vypočteme jako

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Libovolný korelační koeficient nabývá hodnot mezi -1 a 1 a platí, že

- absolutní nepřímá závislost má  $\text{Cor}(X, Y) = -1$
- lineární nezávislost/ nekorelovanost má  $\text{Cor}(X, Y) = 0$
- absolutní přímá závislost má  $\text{Cor}(X, Y) = 1$

O statistické významnosti závislosti rozhodujeme testem

- $H_0$  : korelační koeficient = 0
- $H_1$  : korelační koeficient  $\neq 0$ ,  $> 0$ ,  $< 0$

Za platnosti nulové hypotézy platí, že testová statistika

$$T = \frac{\text{Cor}(X, Y)}{\sqrt{1 - \text{Cor}(X, Y)^2}} \sqrt{(n - 2)}$$

má  $t$ -rozdělení o  $n - 2$  stupních volnosti.

# Pearsonův korelační koeficient

V případě, že chceme testovat konkrétní hodnotu korelačního koeficientu, tedy

- $H_0$  : korelační koeficient =  $\rho_0$
- $H_1$  : korelační koeficient  $\neq \rho_0$ ,  $> \rho_0$ ,  $< \rho_0$

pak se využívá tzv. Fisherovy  $Z$ -transformace, která říká, že

$$Z = \frac{1}{2} \ln \left\{ \frac{1 + \text{Cor}(X, Y)}{1 - \text{Cor}(X, Y)} \right\} \sim N \left( \frac{1}{2} \ln \left\{ \frac{1 + \rho}{1 - \rho} \right\}, \frac{1}{n-3} \right)$$

kde  $\rho$  je skutečná/ teoretická hodnota korelačního koeficientu. Pomocí této  $Z$ -transformace je možné porovnávat i dva korelační koeficienty mezi sebou. Platí totiž, že

$$U = \frac{Z_1 - Z_2}{\sqrt{\frac{1}{n_1-3} + \frac{1}{n_2-3}}}$$

má při shodě porovnávaných korelačních koeficientů  $N(0, 1)$  rozdělení.

V případě, že náhodná veličina normální rozdělení nemá, respektive, že odchylky od normálního rozdělení jsou takového typu, že nelze použít zvolený výše uvedený test, je potřeba zvolit odpovídající **neparametrickou metodu**.

Neparametrické testy bývají většinou založeny na pořadí naměřených hodnot v uspořádané řadě.

**Příklad.** *Uvažujme naměřené věky otců 30, 28, 36, 38, 28, 26, 29, 37, 25, 50. Data věků rodičů bývají sešikmena a často obsahují odlehlé hodnoty. Přiřadíme-li hodnotám pořadí podle velikosti, získáme řadu 6, 3.5, 7, 9, 3.5, 2, 5, 8, 1, 10. Takto získaná řada není sešikmená a nemá odlehlé hodnoty.*



# Znaménkový test

Test o hodnotě mediánu jednoho výběru. Testujeme

- $H_0$  : medián =  $m_0$
- $H_1$  : medián  $\neq m_0$ ,  $> m_0$ ,  $< m_0$

Pro každé pozorování spočteme rozdíl  $X_i - m_0$  a spočítáme, kolik těchto rozdílů je kladných. Tento součet označme jako  $Z$ . Za platnosti nulové hypotézy má testová statistika  $Z$  binomické rozdělení  $Bi(n, 1/2)$ , kde  $n$  je počet pozorování.

Pro velká  $n$  je možné použít i transformaci

$$U = \frac{2Z - n}{\sqrt{n}}$$

Která má za platnosti  $H_0$   $N(0, 1)$  rozdělení.

**Příklad.** Pokračujme v příkladu s věky otců 30, 28, 36, 38, 28, 26, 29, 37, 25, 50 a testujme hypotézu, že medián věku otců je 33 let, tj. testujeme

- $H_0$  : medián věku otců je 33 let
- $H_1$  : medián věku otců není 33 let

Spočtíme rozdíly  $X_i - m_0$ : -3, -5, 3, 5, -5, -7, -4, 4, -8, 17.

Kladných hodnot je mezi nimi  $Z = 4$ .  $P$ -hodnota testu vychází 0.75, což je hodnota  $> \alpha (= 0.05)$  a  $H_0$  tedy nezamítáme.

Použitím  $U$ -transformace dostaneme  $U = -0.632$  a  $p$ -hodnotu 0.527.

# Wilcoxonův jednovýběrový test

Znaménkový test porovnává pouze počet hodnot ležících pod mediánem a těch, co leží nad ním. Nezohledňuje však vzdálenost od mediánu. To dělá Wilcoxonův test, neboli **Mann-Whitneyův** test. Ten už je založen na pořadích. Testované hypotézy zůstávají stejné.

## Postup testu

- spočítají se rozdíly od testované hodnoty  $X_i - m_0$
- určí se jejich znaménko
- určí se pořadí absolutních hodnot rozdílů
- spočítá se součet těchto pořadí patřících kladným rozdílům
- označme tento součet  $S^+$  a obdobně označme  $S^-$  součet pořadí pro záporné rozdíly, musí platit  $S^+ + S^- = n(n+1)/2$ .

Pro větší  $n$  lze užít transformaci

$$U = \frac{S^+ - \frac{1}{4}n(n+1)}{\sqrt{\frac{1}{24}n(n+1)(2n+1)}}$$

která má za platnosti  $H_0$   $N(0, 1)$  rozdělení.

**Příklad.** Pokračujme v příkladu s věky otců 30, 28, 36, 38, 28, 26, 29, 37, 25, 50 a opět testujeme hypotézu, že medián věku otců je 33 let, tj. testujeme

- $H_0$  : medián věku otců je 33 let
- $H_1$  : medián věku otců není 33 let

Spočtěme rozdíly  $X_i - m_0$ : -3, -5, 3, 5, -5, -7, -4, 4, -8, 17 a jejich absolutním hodnotám přiřaďme pořadí 1.5, 6, 1.5, 6, 6, 8, 3.5, 3.5, 9, 10. Sečtěme kladné (modré) pořadí  $S^+ = 21$  a záporné (červené) pořadí  $S^- = 34$ . Testová statistika vychází  $U = -0.66$  a  $p$ -hodnota  $0,51 > \alpha (= 0.05)$  a  $H_0$  tedy nezamítáme.

# Wilcoxonův párový test

V případě, že chceme porovnat dva závislé výběry, které nesplňují předpoklad normality, používá se párový Wilcoxonův test.

I zde zůstávají testované hypotézy stejné jako u párového t-testu.

V prvním kroku se spočítají rozdíly v rámci párů, tj. pro každé  $X_i, Y_i, i = 1, \dots, n$

$$R_i = X_i - Y_i$$

Pokud tyto rozdíly nemají normální rozdělení, použije se pro ně jednovýběrový Wilcoxonův test.

# Wilcoxonův dvouvýběrový test

Pro porovnání dvou nezávislých výběrů, které nesplňují předpoklad normality, se používá Wilcoxonův dvouvýběrový test. Testujeme

- $H_0$  : střední hodnota  $X$  – střední hodnota  $Y = 0$
- $H_0$  : střední hodnota  $X$  – střední hodnota  $Y \neq 0, < 0$  nebo  $> 0$

Test je založen na pořadích hodnot sdruženého výběru. Postup

- oba výběry spojí do jednoho sdruženého
- sdružený výběr se uspořádá podle velikosti a každé pozorování dostane své pořadí
- pro oba výběry se vypočte součet pořadí a následně i průměrné pořadí
- pokud jsou si průměrná pořadí podobná, výběry se mezi sebou významně neliší

# Wilcoxonův dvouvýběrový test

Technický výpočet: označme  $T_1, T_2$  součet pořadí v prvním, respektive druhém výběru. Dále vypočteme

$$U_1 = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - T_1, U_2 = n_1 n_2 + \frac{n_2(n_2 + 1)}{2} - T_2,$$

kde  $n_1, n_2$  jsou rozsahy jednotlivých výběrů. Přesný test porovnává hodnotu  $\min(U_1, U_2)$  s kritickou hodnotou. Asymptoticky platí, že

$$U_0 = \frac{U_1 - \frac{1}{2}n_1 n_2}{\sqrt{\frac{n_1 n_2}{12}(n_1 + n_2 + 1)}}$$

má za platnosti  $H_0$   $N(0, 1)$  rozdělení.

# Wilcoxonův dvouvýběrový test

**Příklad.** Chceme porovnat výsledky testů studentů v Ústí nad Labem a v Liberci. Studenti v Ústí dostali bodová ohodnocení 45, 79, 81, 56, 53, 77. Studenti v Liberci získali ohodnocení 76, 62, 84, 80, 41, 79, 66. Testujeme

- $H_0$  : Studenti v Ústí a v Liberci jsou stejní
- $H_1$  : Studenti v Ústí a v Liberci se liší.
- V prvním kroku srovnám všechny hodnoty do řady  
41, 45, 53, 56, 62, 66, 76, 77, 79, 79, 80, 81, 84
- následně jim přiřadím pořadí  
1, 2, 3, 4, 5, 6, 7, 8, 9.5, 9.5, 11, 12, 13
- pak vypočtu  $T_1 = 38.5$ ,  $T_2 = 52.5$ ,  $U_1 = 24.5$ ,  $U_2 = 17.5$ ,  $U_0 = 0.5$ ,  $p = 0.6678$

$P$ -hodnota  $> \alpha$  a tedy nezamítám nulovou hypotézu, neprokázal se rozdíl mezi studenty v Ústí a v Liberci.



# Kruskal-Wallisův test

V případě, že není splněn předpoklad normality při porovnání více než dvou nezávislých výběrů, používá se

**Kruskal-Wallisova ANOVA.** Kruskal-Wallisova ANOVA je přímým zobecněním Wilcoxonova dvouvýběrového testu.

Testujeme

- $H_0$  : Střední hodnoty výběrů se neliší
- $H_1$  : Střední hodnoty výběrů se liší

Stejně jako u dvouvýběrového Wilcoxonova testu srovnáme všechny naměřené hodnoty do řady, určíme jejich pořadí a spočteme statistiky  $T_1, \dots, T_k$ , kde  $k$  je počet výběrů. Pak platí, že testová statistika

$$Q = \frac{12}{n(n+1)} \sum_{i=1}^k \frac{T_i}{n_i} - 3(n+1)$$

má za platnosti  $H_0$   $\chi^2$ -rozdělení.

# Dunnův test

V případě, že Kruskal-Wallisova ANOVA určí, že se výběry mezi sebou významně liší, je potřeba zjistit, které konkrétní dvojice výběrů se liší. K tomu může sloužit např. **Dunnův test**.

Testová statistika porovnávající  $i$ -tý a  $j$ -tý výběr je

$$D = \frac{(|\frac{T_i}{n_i} - \frac{T_j}{n_j}|)}{\sqrt{\frac{n(n+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_j}\right)}}$$

V případě, že v datech jsou shodné hodnoty a je tedy třeba dělit pořadí, používá se statistika

$$D = \frac{(|\frac{T_i}{n_i} - \frac{T_j}{n_j}|)}{\sqrt{\frac{n(n+1) - \sum_{l=1}^r (S_l^3 - S_l)}{12} \left(\frac{1}{n_i} + \frac{1}{n_j}\right)}}$$

kde  $S_l$  je počet  $l$ -té shodné hodnoty.

Tato statistika má za platnosti  $H_0$   $N(0, 1)$ -rozdělení. Pro vícenásobné porovnání se pak použijí upravené  $p$ -hodnoty, aby byla udržena celková hladina testu.

Pokud se chystáme porovnávat několik závislých výběrů, používá se **ANOVA pro opakovaná měření**.

Příklady takovéto situace mohou být

- **Ochutnávka jogurtů:** 20 lidí ochutnává a hodnotí každý všech 5 porovnávaných vzorků jogurtu.
- **Měření opakovaná v čase:** chceme hodnotit vývoj pacientova zdravotního stavu v čase. Pro 30 pacientů děláme opakovaná měření játrových testů.

Stále se testují hypotézy

- $H_0$  : Střední hodnoty výběrů se neliší
- $H_1$  : Střední hodnoty výběrů se liší

# ANOVA pro opakovaná měření

Vyhodnocení hypotéz probíhá opět pomocí porovnaná variability mezi výběry, ale tentokrát s variabilitou zbytkovou. Zbytková variabilita se od celkové variability v rámci výběrů liší tím, že je od snížena o variabilitu způsobenou rozdíly mezi jedinci. Konkrétně se tato zbytková variabilita získá následovně

$$\begin{aligned} SST &= \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{..})^2 = \\ &= \sum_{i=1}^k n_i (\bar{X}_{i.} - \bar{X}_{..})^2 + \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{i.})^2 = \\ &= SSA + SSe \\ SSz &= SSe - SSS = SSe - k \sum_{j=1}^{n_j} (\bar{X}_{.j} - \bar{X}_{..})^2 \end{aligned}$$

Test je pak založen na porovnání  $SSA$  a  $SSz$ .

V případě, že porovnáváme závislé výběry, které nemají normální rozdělení, používá se **Friedmanův test**. Myšlenkou testu je, že každý jedinec přiřadí jednotlivým vzorkům pořadí od 1 do  $k$  (pro hodnoty naměřené v čase se určí pořadí v rámci každého jedince) a tato pořadí se pak sečtou a zprůměrují. Označme tyto průměry  $\bar{r}_{.j}$ . Ty jsou pak základem testové statistiky

$$Q = \frac{12n}{k(k+1)} \sum_{j=1}^k \left( \bar{r}_{.j} - \frac{k+1}{2} \right)^2$$

Za platnosti nulové hypotézy má tato statistika  $\chi^2$ -rozdělení o  $k - 1$  stupních volnosti.

# Spearmanův korelační koeficient

Pokud chceme otestovat, zda spolu souvisí dvě číselné proměnné, které nemají normální rozdělení (ale stále se jeví jako spojené), používá se **Spearmanův korelační koeficient**. Stejně jako další neparametrické testy je zaměřen na pořadích.

Postup

- Hodnoty každé proměnné převedu na pořadí.
- Spočítá se Pearsonův korelační koeficient pro tato pořadí.

Spearmanův korelační koeficient měří monotónní vztah dvou veličin. Je tedy obecnější než Pearsonův korelační koeficient, který měřil jen lineární závislost.

# Kendallův korelační koeficient

Pokud chceme zjistit, zda je lineární vztah mezi dvěma uspořádanými kategorickými proměnnými, používá se **Kendallův korelační koeficient** (Kendalovo  $\tau$ ).

Označme dvě porovnávané proměnné  $X$  a  $Y$ . Nyní uvažujme všechny dvojice naměřených hodnot  $X_i, Y_i$  a pokud pro danou dvojici platí, že  $X_i < X_j$  &  $Y_i < Y_j$  nebo  $X_i > X_j$  &  $Y_i > Y_j$ , pak označme tuto dvojici jakou **souhlasnou**, pokud platí  $X_i < X_j$  &  $Y_i > Y_j$  nebo  $X_i > X_j$  &  $Y_i < Y_j$ , označme ji za **nesouhlasnou**.

**Kendalovo**  $\tau$  je založeno na rozdílu počtu souhlasných ( $n_s$ ) a počtu nesouhlasných ( $n_n$ ) dvojic.

Konkrétně je **Kendalovo**  $\tau$  definováno jako

$$\tau = \frac{n_s - n_n}{n} = \frac{2}{n(n-1)} \sum_{i < j} \text{sign}(X_i - X_j) \text{sign}(Y_i - Y_j)$$

Rozptyl tohoto koeficientu je

$$\text{Var}(\tau) = \frac{2(2n + 5)}{9n(n-1)}$$

a testová statistika  $\tau/\text{Var}(\tau)$  má za platnosti nulové hypotézy asymptoticky  $N(0, 1)$  rozdělení.



Výše uvedený koeficient funguje dobře, pokud v datech nejsou stejné hodnoty. Pokud se stejné hodnoty vyskytnou, používají se následující období tohoto koeficientu.

Pro proměnné se **stejným počtem možných hodnot**

$$\tau_B = \frac{n_s - n_n}{\sqrt{(n_0 - n_1)(n_0 - n_2)}},$$

kde  $n_0 = n(n - 1)/2$ ,  $n_1 = \sum_i t_i(t_i - 1)/2$  a  $t_i$  jsou počty shodných hodnot u proměnné  $X$ ,  $n_2 = \sum_i u_i(u_i - 1)/2$  a  $u_i$  jsou počty shodných hodnot u proměnné  $Y$ .

Pro proměnné s **různým počtem možných hodnot**

$$\tau_C = \frac{2(n_s - n_n)}{n^2 \frac{m-1}{m}},$$

kde  $m$  je minimální počet hodnot u obou proměnných.

Výpočet rozptylů a následných testových statistik pro  $\tau_B$  a  $\tau_C$  je složitý. Přenechme ho tedy softwarům.

Vztah mezi dvěma spojitými proměnnými lze hodnotit i z pohledu **lineární regrese**, která zkoumá příčinnou závislost. V tomto případě máme

- **nezávisle proměnnou**  $X$  – příčinu
- **závisle proměnnou**  $Y$  – důsledek

Předpokládáme lineární model ve tvaru

$$Y_i = \beta_0 + \beta_1 X_i + e_i$$

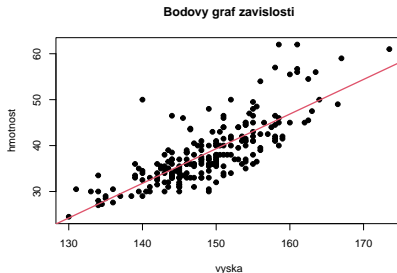
kde

- $Y_i$  jsou hodnoty závisle proměnné
- $X_i$  jsou hodnoty nezávisle proměnné
- $\beta_0$  je absolutní člen
- $\beta_1$  je lineární člen
- $e_i$  jsou náhodné chyby

# Lineární regrese

Graficky popisujeme pomocí bodového grafu, ale není jedno, která proměnná je na které ose

- na x-ovou osu se kreslí nezávisle proměnná
- na y-ovou osu se kreslí závisle proměnná



Odhad probíhá **metodou nejmenších čtverců**, která minimalizuje součet druhých mocnin residuí

$$\min \sum_{i=1}^n R_i^2 = \min \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \min_{b_0, b_1} \sum_{i=1}^n (Y_i - (b_0 + b_1 X_i))^2$$

Hodnoty  $\hat{Y}_i$  se nazývají odhady, nebo též predikce. Hodnoty  $b_0, b_1$  jsou pak odhady regresních koeficientů. Pomocí modelu je možné predikovat budoucí hodnoty závisle proměnné.

Pro hodnotu  $x_0$  nezávisle proměnné  $X$  očekáváme hodnotu

$$\hat{Y}_0 = b_0 + b_1 x_0$$

např. ze známé výšky můžeme predikovat očekávanou hmotnost.

## Koeficient determinace

Zajímavý ukazatel je koeficient determinace

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \text{cor}(X, Y)^2$$

Říká, kolik procent variability závisle proměnné se modelem vysvětlí. Jinými slovy, z kolika procent závisle proměnná závisí na  $X$  a z kolika na něčem jiném.

Na základě modelu lze též zkonstruovat **test nezávislosti**. Testujeme

- $H_0$  : Proměnná  $Y$  (váha) na proměnné  $X$  (výšce) lineárně nezávisí,  $\beta_1 = 0$
- $H_1$  : Proměnná  $Y$  (váha) na proměnné  $X$  (výšce) lineárně závisí,  $\beta_1 \neq 0$

Test je založen na faktu, že  $b_1/\text{se}(b_1) \sim N(0, 1)$ , kde  $b_1$  je odhad lineárního členu  $\beta_1$  a  $\text{se}(b_1)$  je jeho střední chyba.

**Příklad.** *Počítejme závislost hmotnosti na výšce u jedenáctiletých dětí.*

Odhadli jsme model ve tvaru

$$Y_i = -73.81 + 0.75X_i$$

Střední chyba odhadu lineárního členu vyšla 0.04 a testová statistika tedy 18.76. Tu jsme porovnali s kvantilem t-rozdělení  $t_{220}(1 - 0.975) = 1.97$ . Jelikož je testová statistika větší, tak **zamítáme nulovou hypotézu**. P-hodnota testu vyšla  $< 2.2 * 10^{-16}$ , což je menší než  $\alpha = 0.05$ . Koeficient determinace vyšel 0.6153.

**Závěr:** Můžeme tedy říci, že u mužů s jedním rizikovým faktorem ischemické choroby srdeční hmotnost na výšce závisí. Závislost je přímá a vysvětlí se jí 62% variability závisle proměnné (hmotnosti).

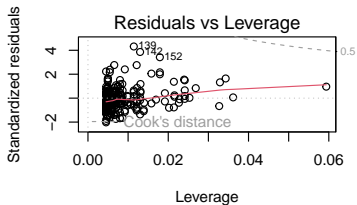
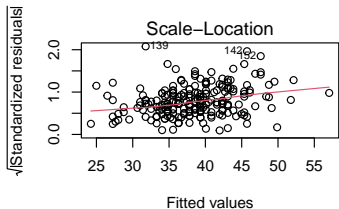
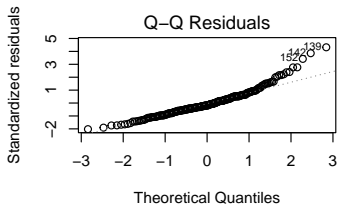
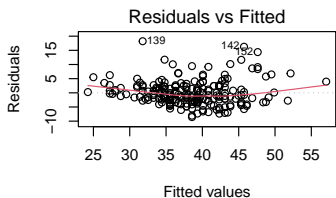
I lineární regrese má své **předpoklady**

- Mezi proměnnými je skutečně lineární vztah
- Residua jsou nezávislá
- Residua mají normální rozdělení
- Stabilita rozptylu
- V datech nejsou vlivná pozorování

Jednotlivé předpoklady můžeme hodnotit buď na základě znalosti dat (nezávislost), nebo grafickými případně číselnými testy.

# Lineární regrese

## Ukázka grafických testů předpokladů





## Ukázka grafických testů předpokladů

- **1. graf:** lineární vztah – červená čára nemá mít trend
- **2. graf:** normalita residuí – body mají ležet na přímce
- **3. graf:** stabilita rozptylu – červená čára nemá mít trend
- **4. graf:** body nemají překročit meze (čárkované křivky)

Regresním modelem nemusíme zkoumat pouze závislost na jedné proměnné, ale můžeme do modelu přidat více nezávisle proměnných. Pak se jedná o **mnohonásobnou lineární regresi** a její model má tvar

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + \epsilon$$

**Příklad.** *Zkoumáme, o kolik stoupne/klesne voda v řece v závislosti na srážkách, na teplotě, na typu půdy, na nasycenosti půdy, na nadmořské výšce, atd.*

Některé proměnné mají na závisle proměnnou větší vliv, jiné menší.

Jak najít optimální model, ve kterém budou jen proměnné s významným vlivem, řeší **kroková regrese**. Tu máme trojího typu

- **Dopředná (forward)**, která začíná s modelem bez nezávisle proměnných a v každém kroku přidá jednu s největším, statisticky významným vlivem
- **Zpětná (backward)**, která začíná s úplným modelem a v každém kroku vynechá jednu proměnnou s nejmenším, statisticky nevýznamným vlivem
- **Kombinace obou předchozích (both sided)**, která začíná s prázdným modelem bez nezávisle proměnných a v každém kroku přidá jednu proměnnou s největším, statisticky významným vlivem a poté zkontroluje, zda nelze jinou proměnnou vynechat.

Cílem je získat model, kde budou nezávisle proměnné pouze se statisticky významným vlivem.

Do modelu mi může vstoupit i kategorická proměnná. Má-li kategorická proměnná  $k$  kategorií, do modelu dáme  $k - 1$  pomocných *dummy* proměnných. Jedná se o 0-1 proměnné a platí

$$\begin{aligned} X_i &= 1 \dots \text{nastala } i\text{-ta kategorie} \\ &= 0 \dots \text{jinak} \end{aligned}$$

Poslední  $k$ -tá kategorie nastane, pokud všechny

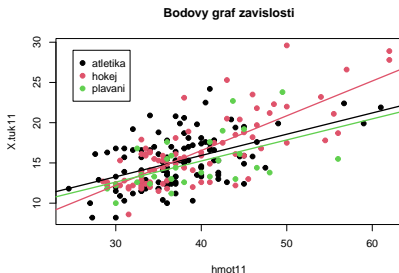
$$X_1 = \dots = X_{k-1} = 0.$$

Každá z těchto *dummy* proměnných má svou  $p$ -hodnotu.

Významnost vlivu kategorické proměnné jako celku se řeší přes **tabulku analýzy rozptylu**, která jí přiřadí jen jednu  $p$ -hodnotu.

Do modelu mohou vstupovat i tzv. **interakce**. Ty popisují způsob, jímž se dvě nezávisle proměnné ovlivňují při jejich současném vlivu na proměnnou závislou.

**Příklad.** *Do výběru bylo zařazeno 222 jedenáctiletých dětí a bylo u nich zjišťováno, jak závisí procento tuku v těle na jejich váze a na sportu, kterému se věnují.*



**Interakce** jsou vidět již z grafu, do kterého se vykreslí závislost číselných proměnných zvlášť pro každou kategorii proměnné kategorické. Pokud interakce v datech jsou, pak se zobrazí různoběžné přímky. Pokud interakce v modelu nejsou, ale skupiny se od sebe liší, zobrazí se rovnoběžné přímky. Pokud rozdíl mezi skupinami není, přímky splývají.

**Příklad.** *V modelu interakce existují - závislost procenta tuku na hmotnosti je jiná pro lední hokejisty a pro ostatní.*

## Informační kritéria hodnotící model.

Každé z níže uvedených kritérií je založeno na věrohodnosti modelu ( $L$  - *likelihood*), tj. na ukazateli, jak dobře model kopíruje data. Tato věrohodnost se dále penalizuje počtem parametrů použitých v modelu  $k$ . Platí, že čím menší hodnota kritéria, tím lepší je model.

- **Akaikeho informační kritérium (AIC):**

$$AIC = 2k - 2 \ln(L)$$

- **Upravené Akaikeho informační kritérium (AICc)** pro malé vzorky:

$$AICc = AIC + \frac{2k(k+1)}{n-k-1}$$

- **Bayesovské informační kritérium (BIC):**

$$BIC = \ln(n)k - 2 \ln(L)$$

Co dělat, když nejsou splněny předpoklady na rozdělení náhodné chyby modelu?

- Závisle proměnná je **spojitá** – použijeme pro závisle proměnnou transformaci, která ji posune k normálnímu rozdělení. Nejčastěji se používá přirozený logaritmus, nebo Box-Coxova transformace.
- Závisle proměnná je **dvouhodnotová** (0-1) – použije se logistická regrese.
- Závisle proměnnou tvoří **počty** – použije se Poissonova regrese.
- Závisle proměnná je **ortogonální** – použije se ordinální regrese.



Závisle proměnná je dvouhodnotová, kdy pravděpodobnost hodnoty 1 označme jako  $\pi$ . A potřebujeme modelovat, jak tato pravděpodobnost hodnoty 1 závisí na dalších ukazatelích.

V tomto případě se odhaduje model ve tvaru

$$\frac{\pi}{1 - \pi} = \exp\{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + \epsilon\}$$

Parametr  $\frac{\pi}{1 - \pi}$  se jmenuje **šance** a počítá se jako pravděpodobnost, že jev nastal, vs. pravděpodobnost, že jev nenastal.

Abychom mohli model odhadnout, modeluje se logaritmus šancí pomocí klasické lineární regrese

$$\ln\left(\frac{\pi}{1 - \pi}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k + \epsilon$$

Pokud z odhadnutého modelu pak chceme zpětně získat vztah pro pravděpodobnost  $\pi$ , použijeme

$$\pi = \frac{\exp\{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k\}}{1 + \exp\{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k\}}$$

Koeficient  $\beta_1$  se interpretuje následovně: "Šance vlastnost mít mi při nárůstu  $X_1$  o 1 vzroste průměrně  $\exp \beta_1$  krát při stejných hodnotách ostatních nezávisle proměnných."

**Příklad.** Uvažujme 150 cestujících na Titaniku. Ke každému cestujícímu máme uvedeno pohlaví, věk, třídu, ve které cestoval a informaci, zda se zachránil nebo ne. Následující tabulky ukazují, jací cestující se zachránili, a jací se utopili.

	Muž	Žena
<i>Přežil</i>	23	26
<i>Nepřežil</i>	89	12

	Dospělý	Dítě
<i>Přežil</i>	46	3
<i>Nepřežil</i>	97	4

	1. třída	2. třída	3. třída	Posádka
<i>Přežil</i>	11	10	11	17
<i>Nepřežil</i>	2	12	39	48

**Příklad.** Označme  $\pi$  pravděpodobnost, že dotyčný přežil.  
Odhad modelu logistické regrese vyšel

$$\ln\left(\frac{\pi}{1-\pi}\right) = 1.45 - 1.58(2.tr) - 3.04(3.tr) - 1.72(posadka) + \\ + 2.32(zena) - 0.9(dospely)$$

*Jako významné vyšly proměnné pohlaví ( $p = 2.55 \times 10^{-6}$ ) a třída ( $p=0.0014$ ) v níž dotyčný cestoval, konkrétně se významně liší třetí třída od první třídy a na hladině významnosti 0.1 i posádka od první třídy.*

*Ženy mají  $\exp(2.32) = 10.17$  krát větší šanci na přežití než muži při ostatních parametrech neměnných.*

*Cestující v první třídě mají  $1 / \exp(-3.04) = 20.9$  krát větší šanci přežít než cestující ve třetí třídě při ostatních parametrech neměnných.*

*Cestující v první třídě mají  $1 / \exp(-1.72) = 5.6$  krát větší šanci přežít než posádka při ostatních parametrech neměnných.*

Většinu základních statistických metod lze zobecnit na mnohorozměrnou situaci.

Předpokládejme, že nemáme jednu proměnnou  $X$ , ale vektor proměnných  $\mathbf{X} = (X_1, X_2, \dots, X_k)^T$ .

**Příklad.** *Měříme několik fyzických parametrů jedince: výška, váha, krevní tlak, vitální kapacitu plic, atd. Každý žák na vysvědčení dostane známku z několika předmětů: čeština, matematika, zeměpis, přírodopis, atd.*

- Namísto jedné střední hodnoty  $\mu$  a jednoho rozptylu  $\sigma^2$  máme vektor středních hodnot  $\mu = (\mu_1, \dots, \mu_k)^T$  a varianční matici  $\Sigma = (\sigma_{ij})$
- odhadujeme je pomocí vektoru průměrů  $\bar{\mathbf{X}} = (\bar{X}_1, \dots, \bar{X}_k)^T$  a maticí  $\mathbf{S} = (s_{ij})$ , kde  $s_{ij} = \text{cov}(X_i, X_j)$  pro  $i \neq j$  a  $s_{ii} = \text{Var}(X_i)$

Zobecnění základních statistických metod.

- Dvouvýběrový test  $\Rightarrow$  **Hotellingův test**
- Analýza rozptylu (ANOVA)  $\Rightarrow$  **MANOVA**
- Korelační koeficient  $\Rightarrow$  **Kanonické korelace**
- Lineární regrese  $\Rightarrow$  **Mnohorozměrná lineární regrese**, kde závisle proměnná má více složek.

Porovnávám střední hodnotu náhodného vektoru ve dvou populacích. Předpokládám nezávislá měření. Testuji

- $H_0$  : vektory středních hodnot se rovnají
- $H_1$  : vektory středních hodnot se nerovnají

Testová statistika má tvar

$$T^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^T \Sigma^{-1} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})$$
$$\Sigma = \frac{(n_1 - 1)\Sigma_1 + (n_2 - 1)\Sigma_2}{n_1 + n_2 - 2}$$

Testová statistika má za platnosti  $H_0$  Hotellingovo  $T^2$ -rozdělení s  $k$  a  $n_1 + n_2 - 2$  stupni volnosti. Toto lze převést na  $F$ -rozdělení.

Obdobně lze zkonstruovat i testovou statistiku pro jednovýběrový test.

Při srovnání více nezávislých výběrů se opět testují hypotézy

- $H_0$  : vektory středních hodnot se rovnají
- $H_1$  : vektory středních hodnot se nerovnají

Stejně jako u jednorozměrné analýzy rozptylu, i ve vícerozměrné verzi je vyhodnocení hypotéz založeno na porovnání variability vysvětlené a nevysvětlené. Existuje několik testových statistik, kde všechny pracují s maticemi

$$\mathbf{W} = \sum_{i=1}^p \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)^T (\mathbf{Y}_{ij} - \bar{\mathbf{Y}}_i)$$

$$\mathbf{B} = \sum_{i=1}^p n_i (\bar{\mathbf{Y}}_i - \bar{\mathbf{Y}})^T (\bar{\mathbf{Y}}_i - \bar{\mathbf{Y}})$$

kde  $p$  značí počet výběrů a  $\bar{\mathbf{Y}}_i$  průměr  $i$ -tého výběru.



Testové statistiky pro MANOVu.

- **Wilkovo lambda**

$$\Lambda_W = \det \left( \frac{\mathbf{W}}{\mathbf{W} + \mathbf{B}} \right)$$

- **Pillayova stopa**

$$\Lambda_P = \text{tr} \left( \frac{\mathbf{B}}{\mathbf{W} + \mathbf{B}} \right)$$

- **Hotellingovo lambda**

$$\Lambda_H = \text{tr} \left( \frac{\mathbf{B}}{\mathbf{W}} \right)$$

při porovnání dvou výběrů se všechny tyto statistiky smrští na Hotellingův dvouvýběrový test.

# Kanonické korelace

Máme dvě skupiny proměnných  $\mathbf{X}$  a  $\mathbf{Y}$  měřených na stejných jedincích a chceme zjistit, zda mezi těmito skupinami je nějaký vztah, případně jaký.

**Příklad.** *Uvažujme dvě různé skupiny lékařských vyšetření a hodnotíme, zda obě tyto skupiny měří to samé, nebo ne.*

Postup je takový, že nejprve hledáme tzv. kanonické proměnné, pro něž platí

$$K_{11} = \mathbf{a}^T \mathbf{X}, K_{21} = \mathbf{b}^T \mathbf{Y}, \quad \text{cor}(K_{11}, K_{21}) = \max\{\text{cor}(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y})\}$$

Toto je první pár kanonických proměnných. Druhý získáme tak, že je kolmý k prvnímu a opět je pro něj korelace maximální, atd. Takto získáme  $k$  kanonických proměnných, kde  $k$  je minimální počet proměnných v první, respektive ve druhé skupině. Korelace příslušné k párům kanonických proměnných se nazývají kanonické a měří se jimi vztah dvou skupin proměnných.

Máme mnohorozměrná data z několika různých populací a chceme najít nejlepší možný způsob, jak na základě dat rozlišit skupiny mezi sebou. Hledáme postup, jak určit skupinu na základě dat.

**Příklad.** *Uvažujme pacienty s různými nemocemi a mějme ke každému skupinu lékařských testů. Chceme pak najít způsob, jak zařadit pacienta do skupiny jen na základě výsledků testů*

## Postup

- pro každou skupinu spočítáme průměrný vektor
- nového pacienta zařadíme do skupiny, která bude mít svůj průměrný vektor nejbližší k pacientovým výsledkům

Jak dobré je určené rozhodovací pravidlo zjistíme na základě klasifikace, tj. zjištění, kolik jednotek jsme přiřadili správně a kolik chybně.

# Diskriminační analýza

Uvažujme pouze dvě populace s průměry  $\bar{\mathbf{X}}_{1,n}$ ,  $\bar{\mathbf{X}}_{2,n}$ . Vzdálenosti od těchto průměrů měříme Mahalanobisovou vzdáleností

$$D(\mathbf{X}, \bar{\mathbf{Y}}) = \sqrt{(\mathbf{X} - \bar{\mathbf{X}})^T \mathbf{V}^{-1} (\mathbf{X} - \bar{\mathbf{X}})}$$

Platí-li

$$D^2(\mathbf{X}, \bar{\mathbf{X}}_{1,n}) < D^2(\mathbf{X}, \bar{\mathbf{X}}_{2,n}),$$

přičadíme pozorování k první populaci, v opačném případě ke druhé. Aritmetickými operacemi lze získat vektor

$$\mathbf{b} = \mathbf{S}^{-1}(\bar{\mathbf{X}}_{1,n} - \bar{\mathbf{X}}_{2,n}),$$

a rozhodovací pravidlo, že pokud

$$\mathbf{b}^T \mathbf{X} - \mathbf{b}^T \frac{\bar{\mathbf{X}}_{1,n} + \bar{\mathbf{X}}_{2,n}}{2} = \sum_{i=1}^k b_i X_i - b_0 > 0$$

pak pozorování patří do první populace.

K tomuto pravidlu mohou přidat ještě apriorní pravděpodobnosti (třeba relativní četnosti nemocí v populaci.)

$$\mathbf{b}^T \mathbf{X} - \mathbf{b}^T \frac{\bar{\mathbf{X}}_{1,n} + \bar{\mathbf{X}}_{2,n}}{2} + \ln \frac{\pi_1}{\pi_2} > 0.$$

# Shluková analýza – hierarchické metody

Mějme mnohorozměrná data a snažme se v nich najít podobnosti, abychom identifikovali různé skupiny pozorování v datech. Cílem je

- najít optimální počet skupin, tak aby mezi nimi byly rozdíly co možná největší, a v rámci skupiny, aby byly hodnoty co nejpodobnější,
- popsat skupiny tak, aby se mezi nimi dalo rozlišovat

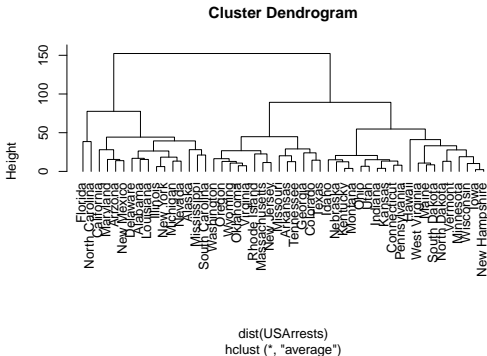
Hierarchické shlukování měří vzdálenosti mezi jednotlivými pozorováními např. euklidovskou vzdáleností a shlukuje k sobě jednotky, co jsou si nejbližší. Vzdálenost skupin se dá měřit trojím způsobem

- vzdálenost středů (průměrů) – **average linkage**
- vzdálenost nejbližších bodů – **single linkage**
- vzdálenost nejvzdálenějších bodů – **complete linkage**

Complete linkage dává většinou nejlepší výsledky.

# Shluková analýza – hierarchické metody

V této analýze nejprve považujeme každé jedno pozorování za samostatnou skupinu a postupně tyto skupiny spojujeme. Graficky se tento proces znázorňuje pomocí **dendrogramu**.



Opticky pak hledáme, kde ukončit shlukování, tj. kolik skupin je optimálních.

Nevýhodou hierarchické metody je, že odlehlé hodnoty v ní často tvoří samostatné skupiny. Alternativou je použít tzv.

**K-means** shlukování. Postup je následující

- nejprve se zvolí počet skupin  $p$
- náhodně vybereme  $p$  bodů v mnohorozměrném prostoru jako středy těchto skupin
- zařadíme prvek, který je nejbližší nějakému středu k této skupině
- středy se přepočítají
- poslední dva body se opakují, dokud nejsou rozřazeny všechny prvky

Nevýhodou tohoto postupu je, že pokud v datech nejsou ednoznačné skupiny, pak rozřazování dopadne jinak při jiné volbě náhodných středů.

Při různých výzkumech bývá často zjišťováno velké množství proměnných, ze kterých má být následně zjištěna nějaká informace. Často bývají mnohé z nich vzájemně korelované a dávají tedy informaci podobnou, ne-li totožnou. Aby bylo možné nějakou informaci z proměnných získat, je často záhodno snížit jejich počet a zabývat se jen těmi skutečně zásadními.



# Metoda hlavních komponent (PCA)

Metoda hlavních komponent transformuje vstupní data tak, aby bylo možné snížit jejich dimenzi / počet. Využívá se přepočítání

$$\mathbf{Y} = \mathbf{X}^T \mathbf{P}$$

kde  $\mathbf{X}$  je centrovaná matice vstupních hodnot (centrování = odečet průměru),  $\mathbf{Y}$  je výstupní - cílová matice a  $\mathbf{P}$  je matice transformačních vektorů. Matici  $\mathbf{P}$  získáme pomocí rozkladu korelační matice vstupních dat  $\mathbf{C}$

$$\mathbf{C} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$$

$\mathbf{\Lambda}$  je pak matice vlastních čísel matice  $\mathbf{C}$  a matice  $\mathbf{P}$  pak obsahuje vlastní vektory matice  $\mathbf{C}$ .

Výsledná matice hlavních komponent  $\mathbf{Y}$  má pak následující vlastnosti

- její vektory jsou vzájemně kolmé (nezávislé)
- řadí se podle variability: od vektoru s největší variabilitou k vektoru s nejnižší variabilitou
- obsahuje veškerou informaci, kterou obsahovala původní data

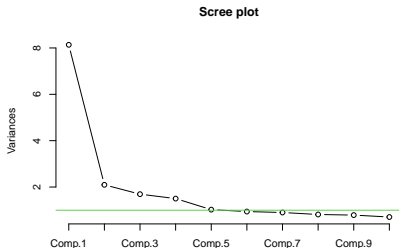
# Metoda hlavních komponent (PCA)

Celý postup si můžeme představit následovně

- představíme si mnohozměrná data v prostoru
- daty proložíme vektor ve směru s největší variabilitou
- tak získáme první hlavní komponentu (PC)
- hledáme vektor, který by byl k prvnímu kolmý a opět byl ve směru s největší variabilitou
- získáme druhou hlavní komponentu
- hledáme vektor, který by byl kolmý k prvním dvěma a byl ve směru s největší variabilitou
- získáme třetí hlavní komponentu
- poslední dva kroky opakujeme, dokud máme body ve volném prostoru

# Metoda hlavních komponent (PCA)

Vstupní data poté reprezentujeme menším množstvím nových proměnných (hlavních komponent) tak, abychom ztratili co nejméně informace / variability. Jejich optimální počet je počet vlastních čísel větších než 1. Graficky znázorněno pomocí tzv. "Scree plot".



Graf zobrazující hodnoty pro prvních 10 hlavních komponent získaných z původních 24 proměnných. Optimální počet hlavních komponent je 5.

Nevýhodou hlavních komponent je, že nemají přirozenou interpretaci. Pokud tedy chceme získat menší počet proměnných, které jsou interpretovatelné, používá se **faktorová analýza**.

Hlavní myšlenka faktorové analýzy pochází z psychologie:

- na každého působí  $k$  neměřitelných faktorů
- podle toho, jak na nás působí, my reagujeme
- podle reakcí na  $p$  podnětů se snažíme identifikovat původní faktory

Vycházíme z rovnice obdobné jako u analýzy hlavních komponent

$$\mathbf{X} = \mathbf{LF}$$

kde  $\mathbf{X}$  je centrovaná matice naměřených dat,  $\mathbf{L}$  jsou tzv. *loadings* a  $\mathbf{F}$  jsou hledané faktory.

Metoda vychází z metody hlavních komponent. Identifikujeme  $k$  hlavních komponent, a ty pak "rotujeme", dokud nedostanou nějakou přirozenou interpretaci. K rotaci je možné použít několik metod, nejčastěji se používá **varimax**.

**Příklad.** Děti nosí ze školy vysvědčení. Podle známek, pak lze identifikovat dvě skupiny studentů, jedna z nich má dobré známky v předmětech *matematika, fyzika, přírodopis, zeměpis, chemie*, druhá má dobré známky v předmětech *čeština, angličtina, dějepis, občanská výchova*. Faktory, které na ně působí jsou pak *přírodní vědy* a *humanitní obory*.

## Popisné statistiky jedné proměnné

- **Spojité proměnná**

- grafy: boxplot, histogram,
- popisné statistiky polohy (průměr, percentily), variability (směrodatná odchylka, mezikvartilové rozpětí, variační koeficient), tvaru rozdělení (šikmost, špičatost)

- **Kategorická proměnná**

- grafy: sloupcový graf, koláčový graf,
- absolutní a relativní četnosti

## Vztah dvou proměnných

- **Dvě kategorické**
  - sloupcový graf
  - kontingenční tabulka, chí-kvadrát test, fisherův test, poměr šancí
- **Spojitá vs. kategorická**
  - boxplot po skupinách, graf průměrů
  - dvouvýběrový test, ANOVA
- **Dvě spojitě**
  - bodový graf
  - korelační koeficient, jednoduchá lineární regrese



## Dvouvýběrové testy – dva nezávislé výběry

- normální data
  - shodné rozptyly – dvouvýběrový **t-test** pro shodné rozptyly (`t.test(..., var.eq=TRUE)`)
  - různé rozptyly – dvouvýběrový **Welchův t-test** pro různé rozptyly (`t.test(...)`)
- nenormální data – **Wilcoxonův** test (je pro něj třeba otestovat shodu rozptylů) (`wilcox.test(...)`)

## **ANOVA** (analýza rozptylu) – více nezávislých výběrů

- normální data
  - shodné rozptyly – klasická **ANOVA** pro shodné rozptyly (`anova(aov(...))`)
  - různé rozptyly – **ANOVA** pro různé rozptyly (`oneway.test(...)`)
- nenormální data – **Kruskal-Wallisův** test (je pro něj třeba otestovat shodu rozptylů) (`kruskal.test(...)`)

## Závislé výběry

- dva výběry
  - normální data – párový **t-test** (`t.test(..., paired=TRUE)`)
  - nenormální data – párový **Wilcoxonův** (`wilcox.test(..., paired=TRUE)`)
- více výběrů
  - normální data – **ANOVA** pro opakovaná měření (`get_anova_table(anova_test(...))`)
  - nenormální data – **Friedmanův** test (`friedman.test(...)`)

## Korelační koeficient – dvě číselné proměnné

- spojitá normální data – **Pearsonův** korelační koeficient (`cor.test(...)`)
- spojitá nenormální data – **Spearmanův** korelační koeficient (`cor.test(..., method="spearman")`)
- kategorická uspořádaná data – **Kendallův** korelační koeficient (`cor.test(..., method="kendal")`)